1 SUPPLEMENTARY INFORMATION

2

# Genomes of all known members of a *Plasmodium* subgenus reveal paths to virulent human malaria

5 **Authors:** Thomas D. Otto[1,†,#,*], Aude Gilabert[2,†], Thomas Crellen[1,3], Ulrike Böhme[1], Céline
6 Arnathau[2], Mandy Sanders[1], Samuel O. Oyola[1,4], Alain Prince Okouga[5], Larson Boundenga[5], Eric
7 Willaume[6], Barthélémy Ngoubangoye[5], Nancy Diamella Moukodoum[5], Christophe Paupy[2], Patrick
8 Durand[2], Virginie Rougeron[2,5], Benjamin Ollomo[5], François Renaud[2], Chris Newbold[1,7], Matthew
9 Berriman[1,*] & Franck Prugnolle[2,5,*]

10 **Affiliations:**

11 [1] Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton CB10 1SA, United Kingdom

12 [2] Laboratoire MIVEGEC, UMR 5290-224 CNRS 5290-IRD 224-UM, Montpellier, France

13 [3] Department of Infectious Disease Epidemiology, Imperial College London, St Mary's Campus,
14 Norfolk Place, London W2 1PG, United Kingdom

15 [4] International Livestock Research Institute, Box 30709, Nairobi, Kenya (current address)

16 [5] Centre International de Recherches Médicales de Franceville, Franceville, Gabon

17 [6] Sodepal, Parc of la Lékédi, Bakoumba, Gabon

18 [7] Weatherall Institute of Molecular Medicine, University of Oxford, John Radcliffe Hospital, Oxford
19 OX3 9DS, United Kingdom

20 # Current Address: Centre of Immunobiology, Institute of Infection, Immunity & Inflammation,
21 College of Medical, Veterinary and Life Sciences, University of Glasgow, Glasgow, United Kingdom

22

23 *Correspondence to: Thomas D. Otto (ThomasDan.Otto@glasgow.ac.uk), Matthew Berriman
24 (mb4@sanger.ac.uk) or Franck Prugnolle (franck.prugnolle@ird.fr)

25

# Contents

49

50

51

# Supplementary Note 1: Dating and population size estimates

A major focus of this study has been to understand the population history of the *Laverania* species and in particular the timing of the divergence events as well as the variation of population size. We used two methods: (1) a Bayesian coalescence model, G-PhoCS[1] to estimate the timing of species divergence (Fig. 1a) and (2) the Multiple Sequentially Markovian Coalescent (MSMC)[2] method to provide a high resolution estimate of changes of $N_e$ through time, specifically to look for a bottleneck that would explain the low diversity in the *P. falciparum* population (Fig. 1b). To scale the population genetic parameters inferred from these models to 'real time' and $N_e$, we used a per-base mutation rate of $3.78 \times 10^{-10}$ (for 4 mitotic events in the red blood cycle)[3].

## Estimation of single nucleotide substitution per year for different generation times

The total number of mitoses per generation was calculated based on different assumptions about total generation time (time to complete the full life cycle), time to complete different stages of the life cycle and number of mitoses per stage.

*Data from previous studies:*

- Development in mosquito takes 10–22 days and involves 10–12 mitoses[4]
- Development in liver takes 5–7 days and involves 15 mitoses[5]
- Gametocyte development takes 12 days and involves 3 mitoses[6].
- Intra-erythrocytic development involves 2 mitoses per day (undergoes three to four rounds of DNA synthesis, mitosis, and nuclear division to produce a syncytial schizont with 16 to 22 nuclei)[6]

74        We have assumed that generation times can be within the range 60 –180 days[7,8].

75    *1.   Assuming 60-day generation time*

| | Days | Min. mitoses | Max. mitoses | Days | Min. mitoses | Max. mitoses |
|---|---|---|---|---|---|---|
| Oocyst to salivary gland | 10 | 10 | 12 | 22 | 10 | 12 |
| Liver | 5 | 15 | 15 | 7 | 15 | 15 |
| Gametocytes | 12 | 3 | 3 | 12 | 3 | 3 |
| subtotal | 27 | 28 | 30 | 41 | 28 | 30 |
| Inferred data, based on 60-day generation time: | | | | | | |
| Blood parameters | 33 | 66 | 66 | 19 | 38 | 38 |
| Total mitoses per gen. | | 94 | 96 | | 66 | 68 |
| Generations per year | | 6.1 | 6.1 | | 6.1 | 6.1 |
| Total mitoses per year | | 572 | 584 | | 401 | 414 |

76    *2.   Assuming 180-day generation time*

| | Days | Min. mitoses | Max. mitoses | Days | Min. mitoses | Max. mitoses |
|---|---|---|---|---|---|---|
| Oocyst to salivary gland | 10 | 10 | 12 | 22 | 10 | 12 |
| Liver | 5 | 15 | 15 | 7 | 15 | 15 |
| Gametocytes | 12 | 3 | 3 | 12 | 3 | 3 |
| subtotal | 27 | 28 | 30 | 41 | 28 | 30 |
| Inferred data, based on 180-day genome time: | | | | | | |
| Blood parameters | 153 | 306 | 306 | 139 | 278 | 278 |
| Total mitoses per gen. | | 334 | 336 | | 306 | 308 |
| Generations per year | | 2 | 2 | | 2 | 2 |
| Total mitoses per year | | 677 | 681 | | 621 | 625 |

77

78   Picking extreme values from **1** and **2** (in red), *total mitoses per year* = 401 to 681

79

80   Using data from Claessens *et al*[3]:

81   Average mutation rate = **3.83 x10$^{-10}$** per base **per 48 hr cycle**

82   (equivalent to 1.64 mutations per genome per year *in vitro*)

83   =>Average mutation rate = **9.57 x10$^{-11}$** per base **per mitosis**

84

85   =>Expected **number of mutations** = (9.57 x10$^{-11}$ x **401**) to (9.57 x 10$^{-11}$ x **681**)

86   = 3.84 x 10$^{-8}$ to 6.52 x 10$^{-8}$ **per base per year**

87   = 0.9 – 1.5 **per genome per year** (considering a genome size of 23.3 Mb)

88  According to Bopp et al[9], excluding parasites grown in presence of drug, the numbers of measured

89  mutations per genome per year were 5.046, 1.682 and 1.682 depending on the isolate. The median

90  value from this study is also nearly identical to that described by Claessens *et al*[3].

## In-vivo data

92  In the *Plasmodium falciparum* IT[10] genome, we observed a region of around 225 – 312 kb, covering

93  the PfCRT locus and an internal *var* gene cluster that is highly conserved in a number of field isolates.

94  Since all of these isolates have the chloroquine resistant genotype, the conserved region is likely to

95  have resulted from chloroquine-selective sweep and could be around 50 years old[11]. However, the

96  presence of a *var* gene on the opposite strand differentiates these isolates from others and may have

97  decreased overall recombination rates in this region.

98  We called SNPs in this region from 5 isolates produced using PacBio sequencing data from the

99  unpublished Pf3k project, available from:

100  ftp://ftp.sanger.ac.uk/pub/project/pathogens/Plasmodium/falciparum/PF3K/PilotReferenceGenomes/)

101  The detected region that is almost SNP free is shown in Supplementary Fig. 2 and the table below. We

102  observed between 0-10 substitutions per year, with a median of 1.67 mutations per genome per year.

103  SNPs were called with *mpileup* and *varfilter* from samtools[12], after remapping the reads with BWA[13].

## Conserved regions around the PfCRT in five clinical *P. falciparum* isolates.

105  Assuming that the selection occurred ~50 years ago, we obtain the reported estimate of mutations per

106  year and genome.

| Isolate | Country | Location of valley on PfIT_07 | Size of valley | mpileup SNPs | Manually inspected SNPs | estimated mutation rate / year / genome |
|---|---|---|---|---|---|---|
| SenT128.08 | Senegal | 383941..609688 | 225 | 0 | 0 | 0 |
| PA0085-C | The Gambia (1) | 378352..699841 | 312 | 1 | 1 | 1.57 |
| PM0138-C | Mali | 372520..664849 | 292 | 1 | 1 | 1.68 |
| PA0012-C | The Gambia (2) | 378563..660686 | 283 | 6 | 6 | 10.38 |
| PD0469-C | Thailand | 228905..510000 | 281 | 4 | 2 | 3.48 |

107  Samples are from the Pf3K pilot project. For further details see
108  ftp://ngs.sanger.ac.uk/production/pf3k/release_5/pf3k_release_5_metadata_20170804.xlsx

109  In conclusion, we assume that *P. falciparum* accumulates on average 0.9-1.5 SNPs per genome per

110  year. We assume that this value is also valid for the other *Laverania* species.

111

## Coalescent models

To estimate key population genetic parameters: effective population size ($N_e$), dates of divergence ($D$) and number of migrants per generation from source population to target population ($M$), we used the Generalised Phylogenetic Coalescent Sampler (G-PhoCS)[1]. As input, we used 1750 alignments from the Lav15sp dataset. We ran two models: (1) split between *P. falciparum* and *P. praefalciparum*, and (2) the entire *Laverania* tree. We incorporated phylogenetic information and modelled bi-directional migration between all extant and ancestral nodes. The MCMC chains were run for a minimum of 10 million iterations, with 20 chains run in parallel. The chains were merged and manually checked for convergence (Tracer version 1.5). We estimated $N_e = \vartheta \,/\, 2\mu$, $D = g^*\tau \,/\, \mu$ and $M = m_{st} * \tau$, where $\vartheta$, $\tau$ and $m_{st}$ *(migration source to target)* are model parameters, $\mu$ is the mutation rate per base pair per generation (ranges from $6.952 \times 10^{-9}$ to $1.158 \times 10^{-9}$ per base-pair per generation, equivalent to 0.9 - 1.5 mutations per year per genome) and $g$ is the generation time of 0.18 to 0.5, as described above. The $M$ parameter is estimated as the total migration rate, approximately indicating the probability that a given lineage in the source population will migrate into the target population[14]. This migration can be seen in some cases (see Supplementary Table 3) especially from *P. praefalciparum* into *P. falciparum*.

We applied the algorithms to three types of alignments (see Supplementary Table 3): (1) genic regions and (2) intergenic regions with and without assumed 500-bp untranslated regions. These alignments appeared robust for the *P. reichenowi*, *P. praefalciparum* and *P. falciparum* comparison as well as for *P. adleri* and *P. gaboni.* However, alignment of more distantly related species was not possible due to a high number of insertions and deletions and the low GC content. We performed the dating on genic alignments for all possible species for which we had more than 2 samples (thus excluding *P. blacklocki*). For the estimates used in Fig. 1 and Supplementary Table 3, some of the estimates of population genetic parameters were approximated where we were unable to generate intergenic alignments.

Multiple Sequentially Markovian Coalescent

To estimate changes in effective population size (*Ne*) over time in *P. falciparum* (PfGA01 & PfIT, from Pf3K dataset), *P. praefalciparum* (PprfG02 & PprfG01) and the gene-flow between them, we ran the multiple sequentially Markovian coalescent (MSMC) on segregating sites from all chromosomes[2]. Genome-wide SNPs were generated by firstly mapping raw reads from each sample against the Pf3D7 reference, then piping BAM files through mpileup v. 0.1.9 (parameters -q 20 -Q 20 -C 50) into bcftools call v. 1.1 (see MSMC documentation for more details). Retaining only homozygous SNPs,

143 each *Plasmodium* chromosome was considered a single phased haplotype. MSMC was run for 20

144 iterations with a fixed recombination rate. Effective population size was calculated as $(1/\lambda)/2\mu$, scaled

145 time was converted into years as *(scaled time / μ) * g*. The parameters λ and scaled time are derived

146 from the model. Values for parameters *μ and g* are described above. The error around our estimates

147 was estimated by bootstrapping 50 replicates by randomly resampling from the segregating sites used

148 as input.

## Estimation of population size

150 Effective population size was estimated from 10,000 years before present (BP) until 500 years BP as

151 bootstrapping demonstrated that the model loses resolution for more recent periods than 500 years BP.

152 The effective population size of *P. falciparum* drops from at least 11,000 years BP and steadily

153 declines to reach its lowest value around 6,000-4,000 years BP ($N_e \sim 3000$), before the population size

154 begins to expand thereafter until 500 years BP (Fig. 2b). While others have speculated on the census

155 population size of *P. falciparum* at this time[15] there is no straightforward way to relate $N_e$ to census

156 population size (*N*) due to complexities in the life-cycle of *P. falciparum* that causes the population to

157 deviate from certain assumptions of the Wright-Fisher model[16]. Nevertheless, generally the census

158 number of parasites is much higher than $N_e$[17]. The bottleneck is unique to *P. falciparum*; although there

159 is a large degree of error in the bootstrapping around $N_e$ estimates for *P. praefalciparum*, the gorilla-

160 infective species does not appear to go through a bottleneck during this period. We replicated the

161 analysis with different *P. falciparum* genomes (PfDd2 & PfHB3, from the Pf3K dataset), which

162 produced near-identical results.

163 Based on evidence from selection in the human genome, the origin of human malaria has been

164 estimated as ~40,000–60,000 years BP and a population expansion associated with the origins of

165 agriculture is assumed to have taken place ~4,000–6,000 years BP[18]. This scenario is confirmed by our

166 modelling of the speciation event between *P. falciparum* and *P. praefalciparum* with G-PhoCS

167 estimates of the timing of speciation ranking from 40,000–60,000 years BP [30,000–70,000 (95 % CI)]

168 and the MSMC estimates of $N_e$ through time showing a rise from 4,000-6,000 years BP onwards.

## Dating of *eba-175* dimorphism

170 To adapt the dating of the *eba-175* dimorphism[19], the following calculation was performed. Previous

171 authors used 6 million years BP as the time when *P. reichenowi* split from the ancestor of

172 *P. falciparum* and *P. praefalciparum* and then dated the *eba-175* split to 0.13–0.14 MYA. As those

173 numbers can be scaled linearly, we used time of 0.13 - 0.23 MYA for the *P. reichenowi* split, which

174    puts the data of the *eba-175* split to around 3,000-5,000 thousand years ago. This agrees with our

175    observation that the dimorphism of *eba-175* occurred in *P. falciparum*, not *P. praefalciparum*

176    (Supplementary Fig. 3b), concluding that the dimorphism occurred during the expansion of the *P.*

177    *falciparum* and its host.

178

179

## Supplementary Note 2: Evolution of core genes

### Within-species polymorphism

The nucleotide diversity per CDS ($\pi$), the average number of nucleotide differences per site between two sequences, was calculated for each species (Supplementary Fig. 1) and their means compared using non-parametric Wilcoxon rank sum tests. Differences in the observed nucleotide diversity may reflect variation in prevalence and different demographic histories of the great ape parasites. The *P. falciparum* nucleotide diversity (computed from 5 worldwide isolates) was significantly lower than the nucleotide diversity observed in any other great ape species (calculated from 2 to 5 genotypes collected in the same localization from Gabon) ($p < 0.0001$). All Wilcoxon rank sum test results comparing nucleotide diversity between the parasites of great apes were significant ($W = 48193000$, $p < 0.0001$; Supplementary Fig. 1). The nucleotide diversity observed in gorilla-infecting species was higher than the diversity observed in the chimpanzee-infecting species ($W_{P.praefal.-P. adleri} = 4963900$, $p < 0.0001$). Among the gorilla-infecting species *P. praefalciparum* presented higher diversity than *P. adleri* ($W_{P. adleri.-P. praefal.} = 9540900$, $p < 0.0001$), due to a higher number of genes with relatively high values of nucleotide diversity. When considering only genes with a nucleotide diversity $\leq 0.02$, the diversity was higher in *P. adleri* ($W_{P. adleri.-P. praefal.} = 9540900$, $p < 0.0001$). Regarding chimpanzee-infecting species, the diversity was significantly higher in *P. gaboni* ($W_{P. reichenowi-P. gaboni} = 5719500$, $p < 0.0001$) and lower in *P. billcollinsi* ($W_{P.reichenowi-P. billcollinsi} = 8011900$, $p < 0.0001$). The lowest diversity was observed in the least prevalent species, *P. billcollinsi* that infects chimpanzees.

### Interspecific gene transfer

Most of the CDS topologies (4,319 out of 4,350, 99.6%, "Lav7sp" dataset) did not significantly differ from the *Laverania* species tree. For the remaining CDS (n=31, including 4 genes of chromosome 4, Supplementary Fig. 5), we specifically looked at their topology and identified those with possible events of gene transfer between species parasitizing the same host species. We detected a clustering of divergent species infecting the same host for eleven CDS, but none of them included all the species infecting the same host, *i.e. P. adleri*, *P. blacklocki* and *P. praefalciparum* or *P. gaboni*, *P. billcollinsi* and *P. reichenowi*. Four of them, localized in the same region of the chromosome 4, shared the same topology, with *P. praefalciparum* and *P. falciparum* grouping together with *P. adleri*, and corresponded to the previously reported introgressed genomic island (topology B in Supplementary Fig. 4; see main text). In the other cases, the chimpanzee-infecting species *P. billcollinsi* was closer to clade A species (four genes; topology C in Supplementary Fig. 5) or clustered together with *P.*

211     *reichenowi* (3 genes; topology D in Supplementary Fig. 5). All these signals remained when

212     considering all sequenced genomes (dataset Lav15st) and concerned in some instances the intergenic

213     region too (see Supplementary Fig. 5, the table below). Beyond these cases, most often, deviations of

214     gene tree topologies from the species tree involved a clustering of *P. billcollinsi* and/or *P. blacklocki*

215     closer to *P. adleri* and/or *P. gaboni* compared to the species phylogeny (not shown), or concerned

216     alignments without enough resolution.

## Genome-wide test of convergent evolution

218     We searched for an excess of convergent substitutions in specific branch-pairs by analyzing the

219     correlation between the number of convergent and divergent substitutions between all the branch-pairs

220     in a phylogeny, and looking for outlier branch-pairs that had high positive residuals, indicating an

221     excess of convergent substitutions relative to the number of divergent substitutions[20]. Both for the

222     divergent and convergent substitutions and for all pairwise comparisons, Pearson's correlation

223     coefficients between the number of substitutions estimated under distinct evolutionary models were

224     always higher than 0.99. We therefore only report the results obtained under the LG model of amino-

225     acid substitutions. At a chromosome scale, we did not detect an excess of convergence between

226     parasite species infecting the gorillas or between the parasites infecting the chimpanzees. However, we

227     detected an excess of convergent substitutions relative to divergent substitutions, in three branch-pairs

228     involving *P. blacklocki* but with no association with the host species.

229

230

## Supplementary Note 3: Gene family analyses

### Differences in gene families

The *P. reichenowi, P. gaboni* and *P. adleri* reference genomes are from single *Laverania* infections where a single isolate predominated (see Supplementary Table 1). However, the *P. praefalciparum* sample contained two distinct genotypes of *P. praefalciparum*. For the core region, a single haploid assembly could be resolved into the two genotypes. For more variable regions of the genome, like the subtelomeres, the genotypes could not be completely resolved and the numbers reported for the *rif, stevor* and *var* genes therefore contain contributions from both haplotypes. For *P. billcollinsi* and especially for *P. blacklocki*, we could not estimate the extent to which the subtelomeres assembled. Although gene families, like CLAG and the *var* genes from internal clusters, did assemble, the numbers of variable genes families are likely to be underestimated due to amplification biases introduced by the sWGA approach for *P. blacklocki* and the fact that *P. billcollinsi* is obtained from a co-infection with *P. gaboni*.

To estimate the number of genes we used (a) a regular expression to count the genes based on functional annotation and (b) matches to Pfam domains (E-value < 1e-6). To each gene/domain we associated counts and standard deviations (Supplementary Tables 6a, b). Differentially distributed gene families are reported in Fig. 3. For several genes, we performed phylogenetic analyses (Supplementary Fig. 7) to better understand their evolution. This was done by aligning the genes of a specific group with Muscle[21] using default parameters. In Seaview[22], we ran GLOCKS[23] with permissive settings and PhyML[24] (default settings for amino acids) to construct trees. The obtained trees were analysed in Figtree[25].

To perform the alignments of *msp* and *eba-175* dimorphic alleles, the same method was used but the numbers of sequences were reduced by subsampling to visualize dimorphisms.

### Generation of similarity matrices

Where sequences were too divergent to perform tree based analyses, we implemented a visualization method based on similarity scores. First, amino acid sequences were compared with a BLASTp (e-value < 1e-6 and low complexity filter set to false). A similarity matrix based on the score or the global identity was built (the alignment length was normalized by the mean sequence length). Using the similarity matrix, the aligned sequences were clustered using the ward.D2 algorithm in the heatmap.2 module of gplots in R[26]. To each gene, we associated their species and in some case their functional

261  annotation through further heatmaps. We used this approach to analyse domains of *var* genes (see
262  below, Supplementary Fig. 10).

## The Rifin and Stevor proteins

264  To build a BLAST-based network, all Pir proteins were compared with an all-against-all BLASTp
265  (parameter: -e 1e-6 –F F). We clustered the Pir proteins using Gephi[27] and tribeMCL[28] into groups,
266  used in Fig. 3. For the Stevor proteins, we built a phylogenetic tree, using RAxML, with the
267  PROTGAMMAIGT model and 100 bootstraps.

### *Meme-Motif analysis for Stevor proteins*

269  To predict motifs in this family, we used MEME[29] version 4.9.1. We searched for 96 motifs of 8-15
270  amino acids using all of the Stevor proteins encoded by the seven reference genomes. Proteins with less
271  than 5 hits were excluded. The output was parsed with a PERL script into a matrix and visualized in
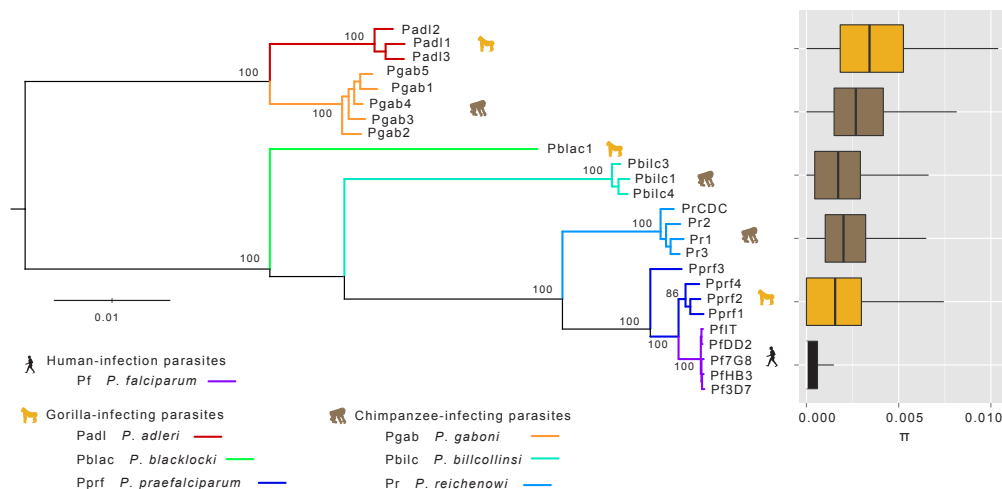272  R[26], using the heatmap.2 function and the ward2 clustering (Supplementary Fig. 8).

## *var* gene analysis

274  To analyse full-length *var* genes in the *Laverania* we excluded genes smaller than 2.5kb and called
275  domains in the genes. The following domains were identified from their conceptual translations: ATS
276  (Acidic Terminal Sequence), NTS (N-Terminal Sequence), DBL (duffy binding like), CIDR (cysteine-
277  rich interdomain region), pam (placenta associated malaria) and the duffy-binding like domain as
278  defined by Pfam (present in invasion related proteins). To call domains, the program hmmscan[30] was
279  used with the HMMer models from the VARdom server using the following parameters: --domT 50 -E
280  1e-6 to attribute domains to *var* genes. As the domains are similar to each other, we generated a PERL
281  program that ascribed domains based on best scores (at least 80% of the length of the HMMer
282  domains). The regions of *var* genes encoding domains could overlap by up to 20 bp. In some cases,
283  rather than finding one of the known domains (DBL, CIDR, ATS or NTS) the Pfam-defined duffy
284  binding-like domain was found. If this happened, we named that domain Duffy, rather than Duffy
285  Binding-Like. Regions (≥ 300aa) in the *var* genes not covered by known domains were also extracted
286  and first called "Unclassified". From those "Unclassified" domains a novel domain was found that we
287  termed CIDRn because of the similarity to existing CIDR domains (Supplementary Fig. 9). To better
288  understand the structure of the domains, particularly Duffy, we used a similarity matrix
289  (Supplementary Fig. 10c).

290    It can be seen that some domains like CIDRα or ATS form defined groups, with little similarity to
291    others. Other domains like DBLα with DBLβ share sequence similarity. The distribution of DBLε,
292    DBLpam2/3 and the unclassified Duffy domain is noteworthy (dotted black lines, top of
293    Supplementary Fig. 10c). These domains seem to be most common in *P. praefalciparum*, *P. adleri* and
294    *P. gaboni*. They have less similarity to other domains. Rather than representing a new domain (like a
295    DBLx[31]), we think that those domains might be more ancient.

296    We also classified the DBLx domain proposed by Larremore *et al*[31]. Their sequences that started with
297    the amino acids specific for the DBLx domain (start NI or DF, end CPQNLDFDRRDQFLR) were
298    compared to our domain dataset. Domains containing those sequences were labelled as DBLx in our
299    set. Next, we generated a similarity matrix with those DBLx, DBLε and Duffy (Supplementary Fig. 12)
300    The DBLx labeled sequences are clustered within the DBLε group. Therefore, we think that the DBLx
301    is not a new domain, but rather part of the diverse DBLε group.
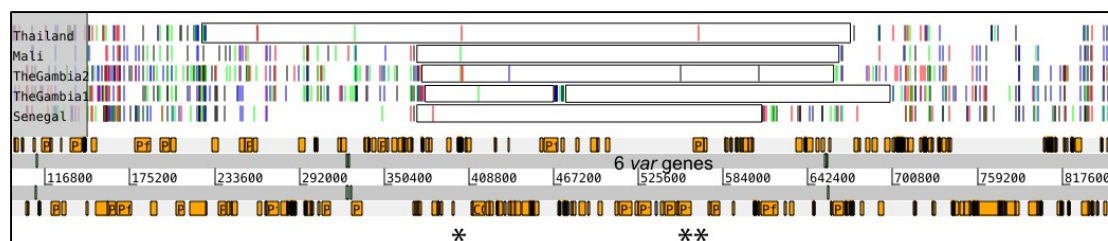
# Supplementary Figures



303

**Supplementary Fig. 1. Maximum Likelihood tree and nucleotide diversity of** ***Laverania*** **isolates.** The tree was obtained using the sequences of 424 genes ("Lav25st" set of orthologues). The box plots show the nucleotide diversity per CDS ($\pi$) for each species. Each boxplot (Tukey's box plot: median, 25th & 75th percentiles and the whiskers extend to the farthest points that are within 1.5 times the interquartile range) is based 3,808 comparisons ("Lav15st" data set).
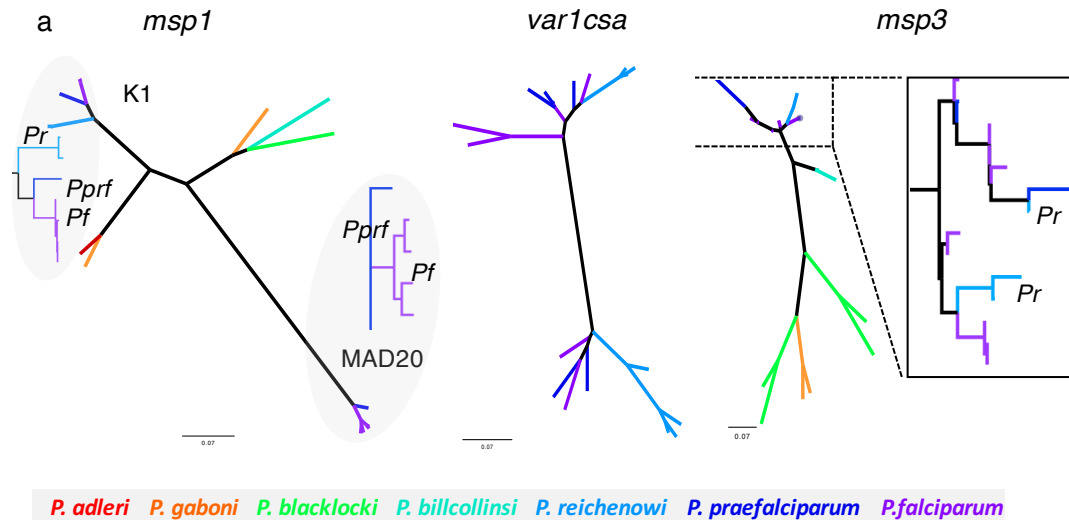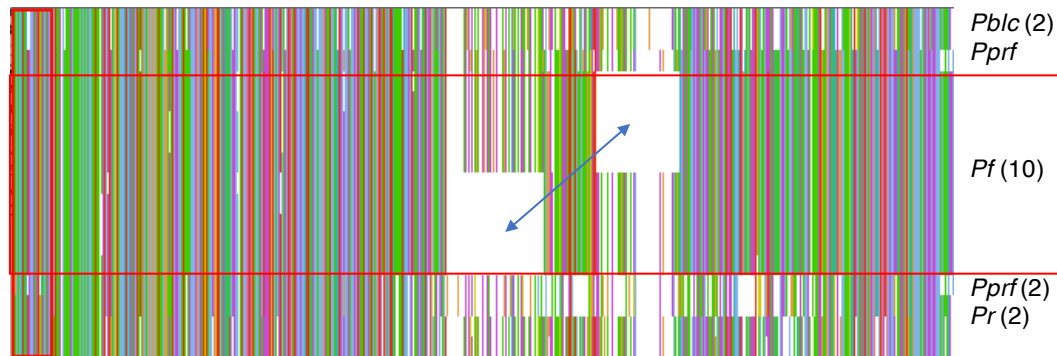
310
311

312

313



**Supplementary Fig. 2. Estimation of** ***in-vivo*** **mutation from field isolates.** An Artemis view of the 600 kb conserved region of five clinical *P. falciparum* isolates (Thailand, Mali, The Gambia 1 and 2 and Senegal), around the *Pf*CRT(*) locus.

14

318    Orange boxes represent the annotated genes on both DNA strands along chromosome

319    7. For each isolate, variation of nucleotide sequences (SNP) compared to the

320    *P. falciparum* 3D7 reference genome are indicated by coloured bars. Large, nearly

321    SNP-free, regions (black boxes) of around 200 kb are found. One *var* gene (**) in the

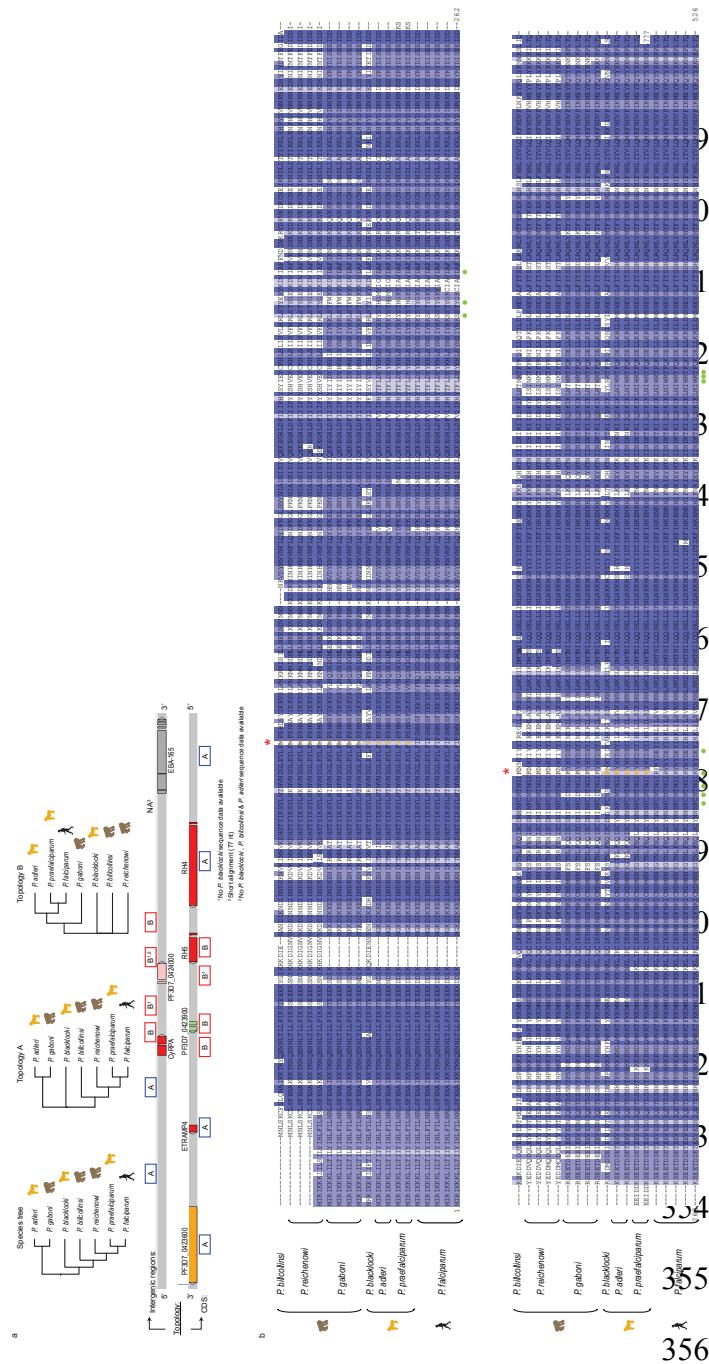322    internal cluster is on the opposite strand.

15

a  *msp1*  var1csa  msp3

P. adleri   P. gaboni   P. blacklocki   P. billcollinsi   P. reichenowi   P. praefalciparum   P.falciparum

(b) *eba-175*

**Supplementary Fig. 3. Dimorphisms in the *Laverania*.** (**a**) Examples of ancient
dimorphisms based on maximum likelihood phylogenetic trees. Dimorphism
in *msp1* arose in the *P. falciparum–P. praefalciparum* ancestor, after the divergence
of *P. reichenowi* and dimorphism in *var1csa* evolved in the *P. reichenowi–P.*
*praefalciparum–P. falciparum* ancestor after the divergence of *P. billcollinsi*. There is
also evidence of a bi-allelic distribution of *msp3* in *P. falciparum, P. praefalciparum*
and *P. reichenowi*. (**b**) Dimorphism in *eba-175* is more recent. The alignment shows
two mutually exclusive indels (arrow) in the *P. falciparum* sequences, not present in
other *Laverania* species. The colours represent different nucleotides. For
the *P. falciparum* sequences, we used full sequences from the following Pf3K
isolates: PfML01, PfSD01, PfDd2, Pf7G8, PfHB3, PfSN01, PfIT, PfCD01, PfGB4

16

336     and PfGN01). *Pf, P. falciparum; Pprf, P. praefalciparum; Pr, P. reichenowi*; and

337     *Pbilc, P. billcollinsi*

338

357

**Supplementary Fig. 4. Interspecific gene transfer and convergent evolution at the right end of the chromosome 4.** (**a**) Support for interspecific gene transfer between the gorilla-infecting species *P. adleri* and the common ancestor of *P. praefalciparum* and *P. falciparum*. The topologies observed in the coding and intergenic regions of the end of chromosome 4 and beyond are given. (**b**) Convergent evolution in the *rh5* gene. Amino acid alignment of the *rh5* region that carries the

18

364   significant fixed difference between parasites infecting the chimpanzees and those

365   infecting gorillas (red stars). Green circles indicate positions that are known to be

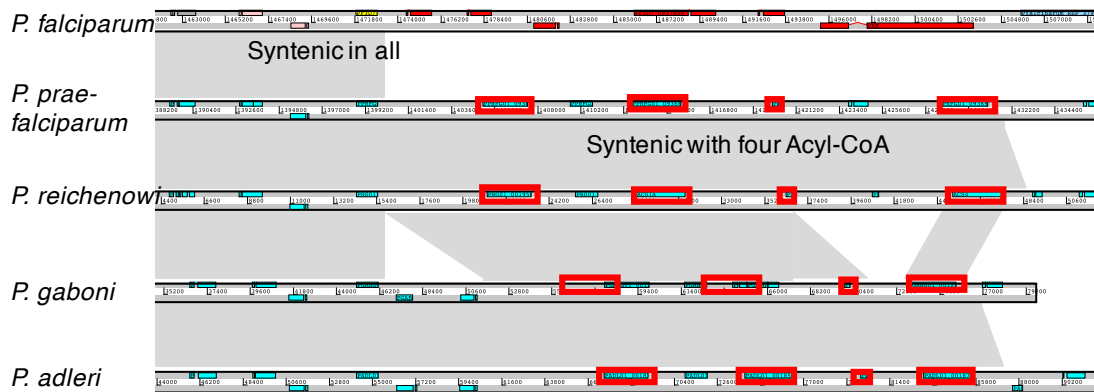366   involved in the interaction with the human receptor Basigin[32].
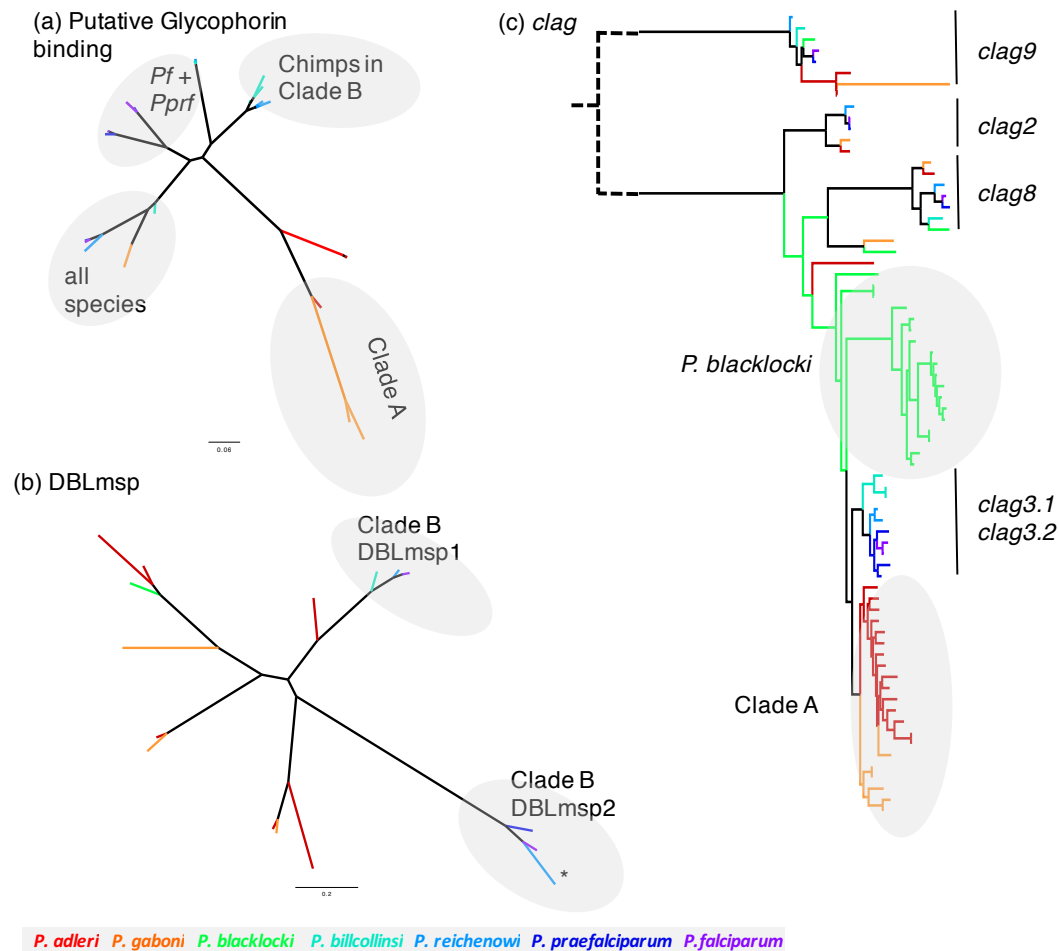
367

368

369

Supplementary Fig. 5 figure: Species tree and Topology A (4319 CDS), Topology B (4 CDS), Topology C (4 CDS), Topology D (3 CDS).

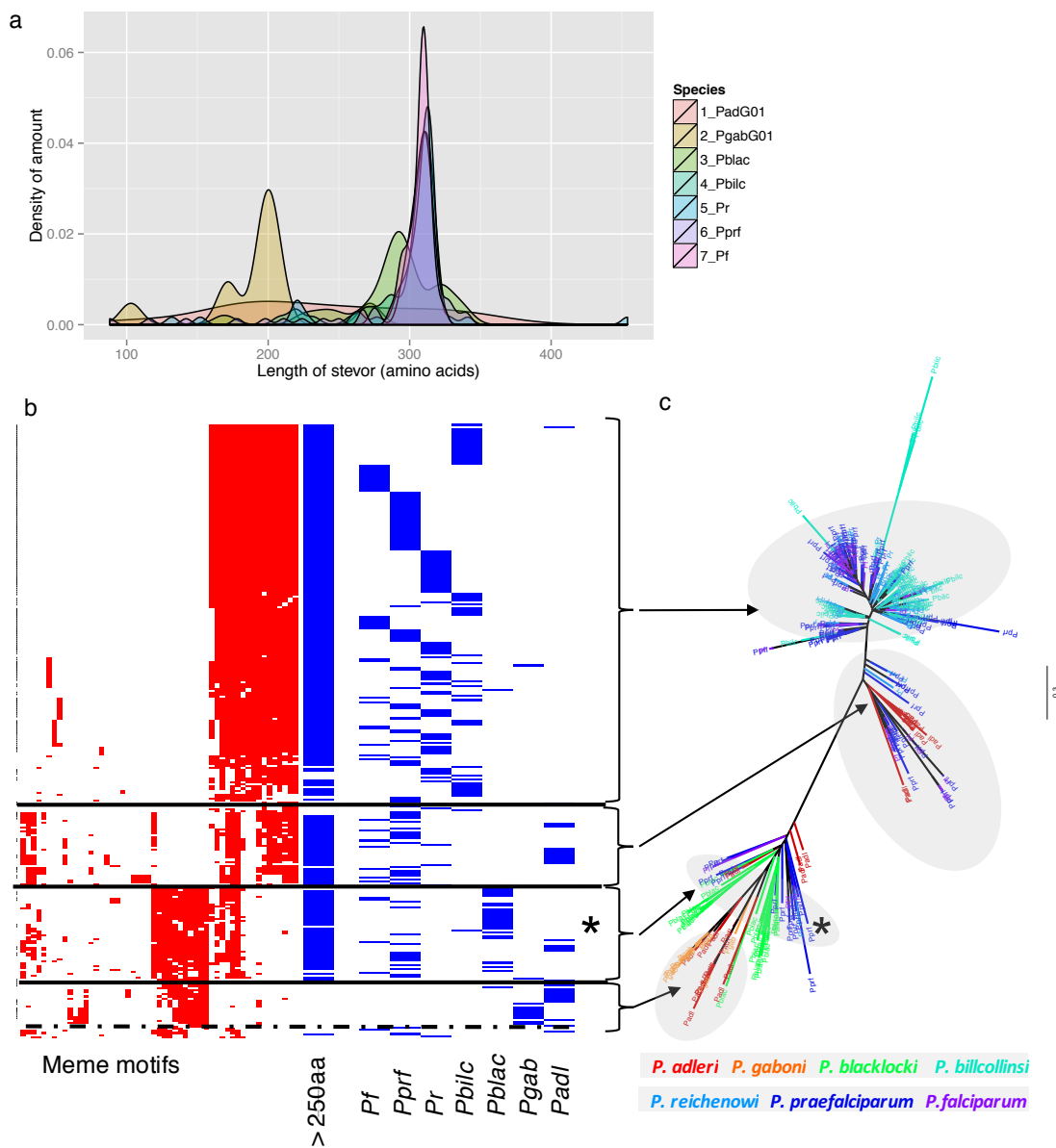| gene ID | function | gene name | topology | newick tree | all strains | | intergenic regions | | note |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | nt | aa | 5' | 3' | |
| PF3D7_0423800 | cysteine-rich protective antigen RH5-Ripr membrane anchoring protein | CyRPA | B | (((PADL01,(PPRFG01,PF3D7)),PGABG01),(PBILCG01,PRG01),PBLACG01); | yes | yes | no | yes | (Pprf, Pf) clusters together with Padl |
| PF3D7_0423900 | probable protein, unknown function | | B | (PBILCG01,((PGABG01,((PF3D7,PPRFG01),PADL01)),PRG01),PBLACG01); | yes | yes | yes | yes | (Pprf, Pf) clusters together with Padl |
| PF3D7_0424000 | Plasmodium exported protein (PHISTc), unknown | | B | ((PGABG01,(PADL01,(PF3D7,PPRFG01))),PBILCG01,PRG01); | yes | yes | yes | yes | No Pblac sequence (Pprf, Pf) clusters together with Padl |
| PF3D7_0424100 | reticulocyte binding protein homologue 5 | RH5 | B | ((((PPRFG01,PF3D7),PADL01),PGABG01),(PRG01,PBILCG01),PBLACG01); | yes | yes | yes | yes | (Pprf, Pf) clusters together with Padl |
| PF3D7_0524000 | karyopherin beta | KASbeta | C | (PPRFG01,(PRG01,(PBLACG01,((PBILCG01,PGABG01),PADL01))),PF3D7); | yes PgabG02 inbetween (Padl,PgabG01,Pbilc) & the others | Pbilc closer to cladeA sp. but outside (Padl,Pgab) | no | NA no Pbilc sequence | Pbilc clusters with Pgab |
| PF3D7_0613200 | conserved Plasmodium protein, unknown function | | C | ((PBLACG01,((PF3D7,PPRFG01),PRG01)),(PGABG01,PBILCG01),PADL01); | yes | NA | yes | yes | Pbilc closer to clade A |
| PF3D7_1327500 | conserved Plasmodium protein, unknown function | | C | ((PGABG01,PADL01),(PBLACG01,(PRG01,(PF3D7,PPRFG01))),PBILCG01); | yes | yes | yes | no | Pbilc closer to clade A |
| PF3D7_1328000 | conserved Plasmodium protein, unknown function | | C | ((((PBLACG01,(PBILCG01,(PADL01,PGABG01))),PRG01),PPRFG01,PF3D7); | yes | yes | no | NA | Pbilc closer to clade A |
| PF3D7_0102400 | lysophospholipase, putative, pseudogene | | D | (PADL01,((PF3D7,PPRFG01),((PBILCG01,PRG01),PBLACG01)),PGABG01); | yes | yes | yes | no | Pbilc & Pr cluster together |
| PF3D7_0102500 | erythrocyte binding antigen-181 | EBA181 | D | ((((PADL01,PGABG01),PBLACG01),(PBILCG01,PRG01)),PF3D7,PPRFG01); | yes | yes | no | no | Pbilc & Pr cluster together |
| PF3D7_0902600 | serine/threonine protein kinase, FIKK family | FIKK9.7 | D | ((((PPRFG01,PF3D7),(PRG01,PBILCG01)),PBLACG01),PADL01,PGABG01); | yes | yes | yes | yes | Pbilc & Pr cluster together |

**Supplementary Fig. 5. Tree topology tests.** The number of protein coding sequences (CDS) producing each of the major tree topologies (observed for >1 CDS) is shown. The table summarizes the results for the 11 CDS with signals of gene flow between species infecting the same host (*i.e.* other signals are not considered here). The table also shows whether a given signal was still observed when all strains were considered, using the nucleotide or amino acid sequences and whether the signal was observed in the intergenic regions down- and up-stream of the respective genes.

**Supplementary Fig. 6. Acyl-CoA Synthetase expansion on Chromosome 9.** ACT view of five genomes, at the right-hand side of chromosome 9. The grey areas indicate co-linearity. *P. falciparum* has lost this region with four Acyl-coA synthetase genes, as this locus is conserved in the other species.

(a) Putative Glycophorin binding

Pf + Pprf

Chimps in Clade B

all species

Clade A

0.06

(b) DBLmsp

Clade B DBLmsp1

Clade B DBLmsp2

*

0.2

(c) *clag*

clag9

clag2

clag8

*P. blacklocki*

clag3.1
clag3.2

Clade A

*P. adleri*   *P. gaboni*   *P. blacklocki*   *P. billcollinsi*   *P. reichenowi*   *P. praefalciparum*   *P.falciparum*

386

387 **Supplementary Fig. 7. Phylogenetic analysis of multigene families.** Example of

388 three families that show differences within the *Laverania*. (**a**) The putative

389 glycophorin binding proteins form four distinct groups. One group contains sequences

390 from all species. The remaining groups are clade, host or species sub-group specific.

391 (**b**) Differences in the DBLmsp that are expanded in Clade A. The DBLmsp2 is a

392 pseudogene (*) in *P. reichenowi*. (**c**) Expansion of *clag* genes in Clade A. The

393 distance between the CLAG9 clade and its nearest neighbour has been compressed to

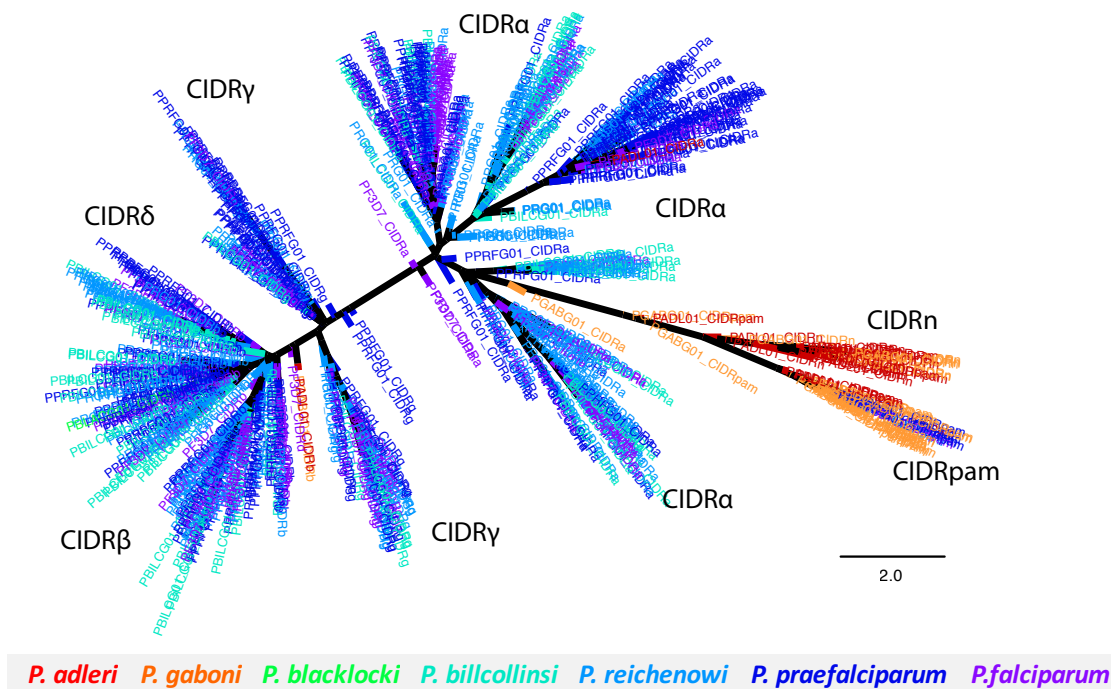394 aid visualisation (dotted lines).

395

22

Supplementary Fig. 8. Analysis of Stevor proteins. (a) Length of Stevor proteins for the seven *Laverania* genomes. (b) Occurrence matrix of meme motifs generated for Stevor proteins (Supplementary note 3). Columns represent the different meme motifs, rows represent all the 301 Stevor proteins. To classify each gene, a binary barcode (blue) is used to indicate whether it encodes likely full length protein (>250aa) and to indicate the species in which it is found. The matrix was clustered with the ward2 algorithm. Note that one cluster (*) has no full length Stevor proteins in chimpanzee parasites. *Pf, P. falciparum; Pprf, P. praefalciparum; Pr, P.*

23

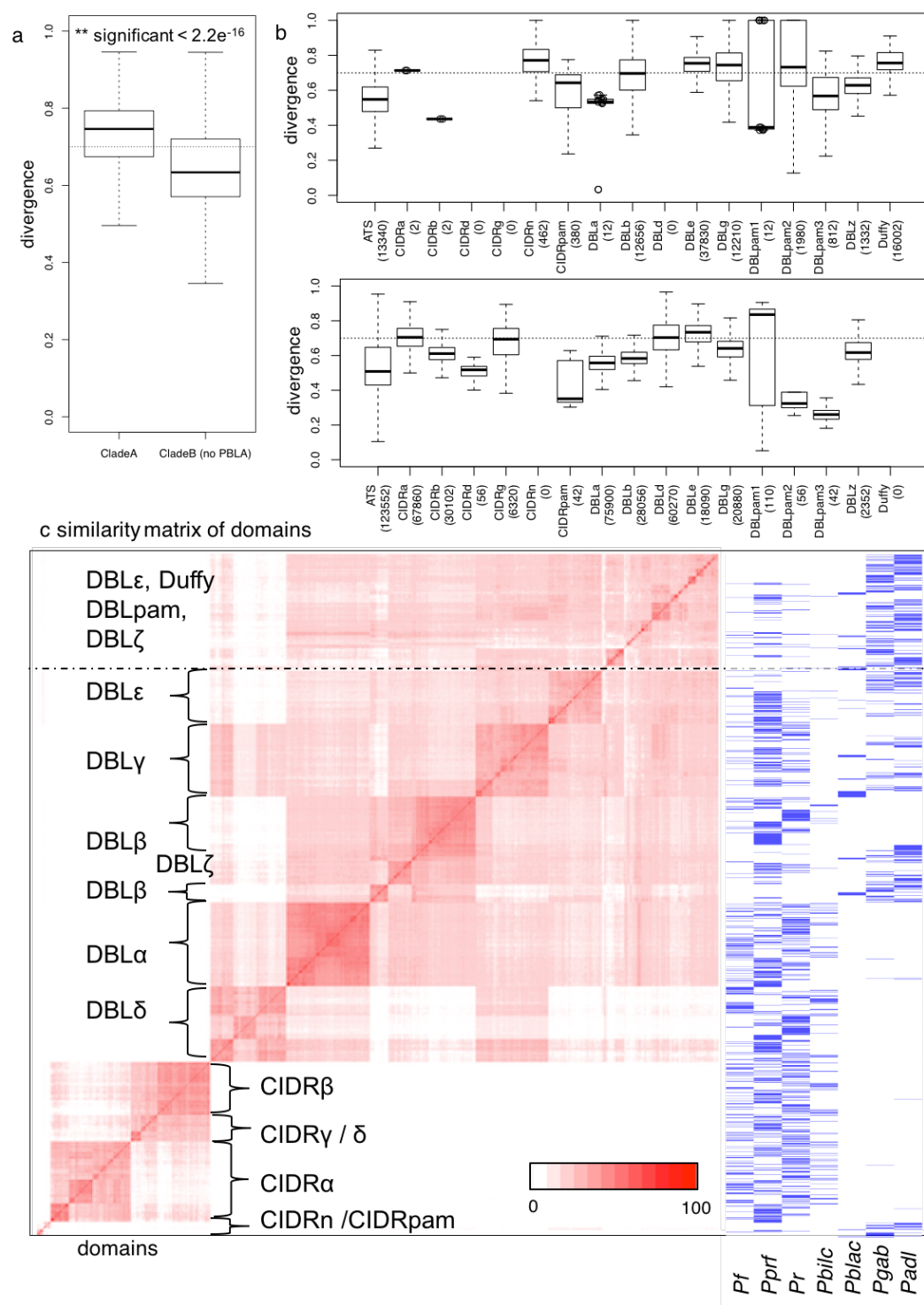405 *reichenowi; Pbilc, P. billcollinsi; Pblac, P. blacklocki; Pgab, P. gaboni;* and *Padl, P.*

406 *adleri*. (**c**) Maximum likelihood tree of the same data. Bootstrap values of 100 were

407 obtained for all branches.

408

409

**Supplementary Fig. 9. Phylogenetic position of the new CIDR domain (CIDRn) specific to Clade A parasites relative to other CIDR domains.** Phylogenetic tree was obtained with RAxML using the PROTGAMMAIGTR models. Bootstrap values of 100 were obtained on all branches.

417

**Supplementary Fig. 10. Diversity of *var* genes domains.** (a) Relative domain
similarity between Clade A and Clade B (excluding *P. blacklocki*) based on the

26

420  average across all domains except ATS. The difference observed between Clade A

421  and Clade B is statistically significant (t-test, two sided). Boxplots are based on

422  83,692 and 310,136 comparisons. (b) Relative similarity across all *var* domain types

423  in Clade A (top) and Clade B (bottom, excluding *P. blacklocki*). The number of

424  predictions for each domain type are shown in parentheses. (c) Annotated similarity

425  matrix between all *var* domains (> 220aa) of the *Laverania* as defined in Fig. 5,

426  including their species and cluster attributions. The similarity matrix shows the score

427  of the BLASTp between the domains, clustered with the ward2 algorithm in R. Each

428  row and column represents one domain and shows its similarity to the other 2,467

429  domains (and itself). The occurrence of each domain/row across the species is

430  indicated by the blue bars on the right. Although domains above the dotted line are

431  classified differently, they cluster together. Pf, *P. falciparum*; Pprf, *P.*

432  *praefalciparum*; Pr, *P. reichenowi*; Pbilc, *P. billcollinsi*; Pblac, *P. blacklocki*; Pgab, *P.*

433  *gaboni*; and Padl, *P. adleri*. All boxplots are Tukey's box plots, showing the median,

434  25th & 75th percentiles and the whiskers and the whiskers extend to the farthest

435  points that are within 1.5 times the interquartile range.

436
437

438

439

440

27

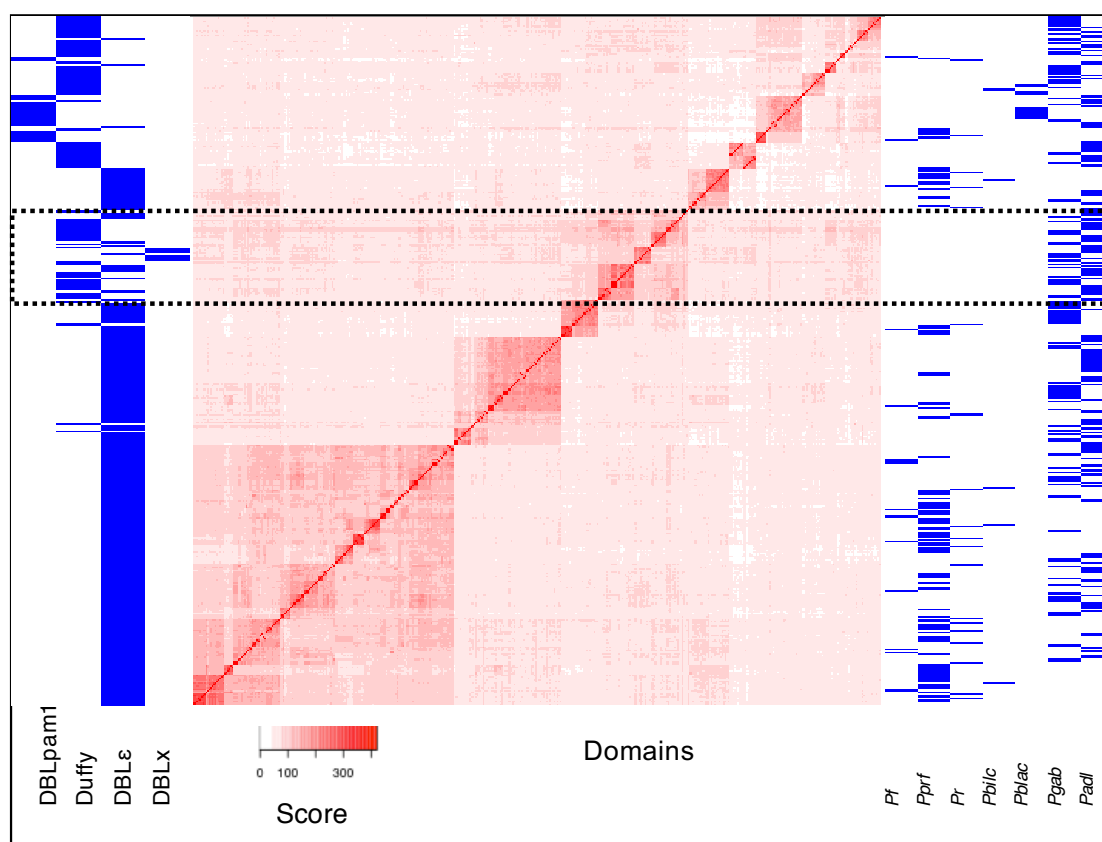**Supplementary Fig. 11. Composition, structure and evolution of *var* genes within *Laverania*.** (a) Screenshot from ACT showing the 2nd internal cluster of *var* genes on chromosome 4 in the seven *Laverania* species. In Clade A and *P. blacklocki*, the orientation of the *var* genes is different compared to that of the other species. The GC-rich RUF elements[33] (RNA of Unknown function), highlighted with an R, occur

28

447     less frequently in Clade A genomes. The size of the *var* genes between the species is

448     different. *var* genes are in red, *pir* genes in green. (b) Bar plot of the number and

449     orientation of *va*r genes or pseudogenes, on the forward (blue) or reverse (red) strand,

450     within internal *var* gene clusters in the *Laverania*. Orientation is relative to the *P.*

451     *falciparum* 3D7 reference genome.

452

453

454

**Supplementary Fig. 12: DBLx is related to DBLε and the ancestral Duffy domain.** To understand the diversity of DBLε and DBLn and to compare them to the newly described DBLx, a similarity matrix of all domains annotated as DBLε, Duffy, DBLpam1 and DBLx labelled sequences, see Supplementary Note 3. It can be seen that all domains are similar to each other and that the DBLx labeled sequences as defined by Larremore *et al*[31] cluster within a group that contains DBLε and Duffy domains (dotted box). *Pf, P. falciparum; Pprf, P. praefalciparum; Pr, P. reichenowi; Pbilc, P. billcollinsi; Pblac, P. blacklocki; Pgab, P. gaboni;* and *Padl, P. adleri*.

455
456
457
458
459
460
461
462

463

464

465

## List of Supplementary Tables 1-9

see associated Excel file.

471      References

472

473    1      Gronau, I., Hubisz, M. J., Gulko, B., Danko, C. G. & Siepel, A. Bayesian inference of ancient
474           human demography from individual genome sequences. *Nat Genet* **43**, 1031-1034,
475           doi:10.1038/ng.937 (2011).
476    2      Schiffels, S. & Durbin, R. Inferring human population size and separation history from
477           multiple genome sequences. *Nat Genet* **46**, 919-925, doi:10.1038/ng.3015 (2014).
478    3      Claessens, A. *et al.* Generation of antigenic diversity in Plasmodium falciparum by structured
479           rearrangement of Var genes during mitosis. *PLoS Genet* **10**, e1004812,
480           doi:10.1371/journal.pgen.1004812 (2014).
481    4      Ponnudurai, T. *et al.* Sporozoite load of mosquitoes infected with Plasmodium falciparum.
482           *Transactions of the Royal Society of Tropical Medicine and Hygiene* **83**, 67-70 (1989).
483    5      Mazier, D. *et al.* Complete development of hepatic stages of Plasmodium falciparum in vitro.
484           *Science* **227**, 440-442 (1985).
485    6      Gerald, N., Mahajan, B. & Kumar, S. Mitosis in the human malaria parasite Plasmodium
486           falciparum. *Eukaryotic cell* **10**, 474-482, doi:10.1128/ec.00314-10 (2011).
487    7      Nkhoma, S. C. *et al.* Population genetic correlates of declining transmission in a human
488           pathogen. *Molecular ecology* **22**, 273-285, doi:10.1111/mec.12099 (2013).
489    8      Molineux, L. in *Malaria, Principles and Practice of Malariology* Vol. 2 (ed McGregor IA
490           Wernsdorfer WH) 913-998 (London Churchill, Livingston, 1998).
491    9      Bopp, S. E. *et al.* Mitotic evolution of Plasmodium falciparum shows a stable core genome
492           but recombination in antigen families. *PLoS Genet* **9**, e1003293,
493           doi:10.1371/journal.pgen.1003293 (2013).
494    10    Udeinya, I. J., Graves, P. M., Carter, R., Aikawa, M. & Miller, L. H. Plasmodium falciparum:
495           effect of time in continuous culture on binding to human endothelial cells and amelanotic
496           melanoma cells. *Experimental parasitology* **56**, 207-214 (1983).
497    11    Payne, D. Spread of chloroquine resistance in Plasmodium falciparum. *Parasitol Today* **3**,
498           241-246 (1987).
499    12    Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-
500           2079, doi:10.1093/bioinformatics/btp352 (2009).
501    13    Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows-Wheeler transform.
502           *Bioinformatics* **26**, 589-595, doi:btp698 [pii] 10.1093/bioinformatics/btp698 (2010).
503    14    Freedman, A. H. *et al.* Genome sequencing highlights the dynamic early history of dogs.
504           *PLoS Genet* **10**, e1004016, doi:10.1371/journal.pgen.1004016 (2014).
505    15    Volkman, S. K. *et al.* Recent origin of Plasmodium falciparum from a single progenitor.
506           *Science* **293**, 482-484, doi:10.1126/science.1059878 (2001).
507    16    Chang, H. H. *et al.* Malaria life cycle intensifies both natural selection and random genetic
508           drift. *Proceedings of the National Academy of Sciences of the United States of America* **110**,
509           20129-20134, doi:10.1073/pnas.1319857110 (2013).
510    17    Palstra, F. P. & Fraser, D. J. Effective/census population size ratio estimation: a compendium
511           and appraisal. *Ecology and evolution* **2**, 2357-2365, doi:10.1002/ece3.329 (2012).
512    18    Karlsson, E. K., Kwiatkowski, D. P. & Sabeti, P. C. Natural selection and infectious disease in
513           human populations. *Nature reviews. Genetics* **15**, 379-393, doi:10.1038/nrg3734 (2014).
514    19    Yasukochi, Y., Naka, I., Patarapotikul, J., Hananantachai, H. & Ohashi, J. Genetic evidence
515           for contribution of human dispersal to the genetic diversity of EBA-175 in Plasmodium
516           falciparum. *Malar J* **14**, 293, doi:10.1186/s12936-015-0820-2 (2015).
517    20    Castoe, T. A. *et al.* Evidence for an ancient adaptive episode of convergent molecular
518           evolution. *Proceedings of the National Academy of Sciences of the United States of America*
519           **106**, 8986-8991, doi:10.1073/pnas.0900233106 (2009).
520    21    Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
521           *Nucleic Acids Res* **32**, 1792-1797, doi:10.1093/nar/gkh340 (2004).

522   22   Gouy, M., Guindon, S. & Gascuel, O. SeaView version 4: A multiplatform graphical user
523          interface for sequence alignment and phylogenetic tree building. *Molecular biology and*
524          *evolution* **27**, 221-224, doi:10.1093/molbev/msp259 (2010).
525   23   Talavera, G. & Castresana, J. Improvement of phylogenies after removing divergent and
526          ambiguously aligned blocks from protein sequence alignments. *Syst Biol* **56**, 564-577,
527          doi:10.1080/10635150701472164 (2007).
528   24   Guindon, S., Delsuc, F., Dufayard, J. F. & Gascuel, O. Estimating maximum likelihood
529          phylogenies with PhyML. *Methods Mol Biol* **537**, 113-137, doi:10.1007/978-1-59745-251-9_6
530          (2009).
531   25   FigTree v.1.4.2, Available http://tree.bio.ed.ac.uk/software/figtree/ (2014).
532   26   Team, R. D. C. *R: A Language and Environment for Statistical Computing*. (2008).
533   27   Bastian, M., Heymann, S. & Jacomy, M. in *International AAAI Conference on Weblogs and*
534          *Social Media*   (2009).
535   28   Enright, A. J., Van Dongen, S. & Ouzounis, C. A. An efficient algorithm for large-scale
536          detection of protein families. *Nucleic Acids Res* **30**, 1575-1584 (2002).
537   29   Bailey, T. L. *et al.* MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res*
538          **37**, W202-208, doi:10.1093/nar/gkp335 (2009).
539   30   Johnson, L. S., Eddy, S. R. & Portugaly, E. Hidden Markov model speed heuristic and
540          iterative HMM search procedure. *BMC Bioinformatics* **11**, 431, doi:10.1186/1471-2105-11-
541          431 (2010).
542   31   Larremore, D. B. *et al.* Ape parasite origins of human malaria virulence genes. *Nature*
543          *communications* **6**, 8368, doi:10.1038/ncomms9368 (2015).
544   32   Wright, K. E. *et al.* Structure of malaria invasion protein RH5 with erythrocyte basigin and
545          blocking antibodies. *Nature* **515**, 427-+, doi:10.1038/nature13715 (2014).
546   33   Guizetti, J., Barcons-Simon, A. & Scherf, A. Trans-acting GC-rich non-coding RNA at var
547          expression site modulates gene counting in malaria parasite. *Nucleic Acids Res* **44**, 9710-9718,
548          doi:10.1093/nar/gkw664 (2016).

549